

AD-A055 236 AIR FORCE INST OF TECH WRIGHT-PATTERSON AFB OHIO SCH--ETC F/G 5/8
AN ALGORITHM FOR DETERMINING SPEECH INTELLIGIBILITY.(U)
DEC 77 W R BEESON

UNCLASSIFIED AFIT/6E/EE/77-9

NL

1 OF 1
AD
A055 236

1



DISCLAIMER NOTICE

**THIS DOCUMENT IS BEST QUALITY
PRACTICABLE. THE COPY FURNISHED
TO DDC CONTAINED A SIGNIFICANT
NUMBER OF PAGES WHICH DO NOT
REPRODUCE LEGIBLY.**

(1)

THIS DOCUMENT IS BEST QUALITY PRACTICABLE
THE COPY FURNISHED TO DDC CONTAINED A
SIGNIFICANT NUMBER OF PAGES WHICH DO NOT
REPRODUCE LEGIBLY.

(6) AN ALGORITHM FOR DETERMINING
SPEECH INTELLIGIBILITY,

THESIS

(9) Master's Thesis,

(14) AFIT/GE/EE/77-9

(10) Wayne R. Beeson
Captain USAF

(11) Dec 77

(16) 7071

(17) $\phi\phi$

(12) 63p.

DDC
RECEIVED
JUN 20 1978
RECEIVED
E

Approved for public release; distribution unlimited

18 06 13 087
012 225

CL

AN ALGORITHM FOR DETERMINING
SPEECH INTELLIGIBILITY

THESIS

Presented to the Faculty of the School of Engineering
of the Air Force Institute of Technology
Air University
in Partial Fulfillment of the
Requirements for the Degree of

Master of Science

ACCESSION for	
NTIS	White Section <input checked="" type="checkbox"/>
DDC	Buff Section <input type="checkbox"/>
UNANNOUNCED	<input type="checkbox"/>
JUSTIFICATION.....	
BY.....	
DISTRIBUTION/AVAILABILITY CODES	
Dist.	AVAIL. and/or SPECIAL
A	23 6

by

Wayne R. Beeson, B.S.

Captain USAF

Graduate Electronic Engineering

December 1977

Preface

This topic was selected because of a continuing need in the R&D community to make an intelligibility evaluation of experimental and prototype voice communications systems. Because the Air Force does not have a more automated way to test intelligibility that produces adequate accuracy and is relatively easy to use, they are still using human listener panels to make these determinations. I have the feeling that "there must be a better way" to make these tests.

It seems reasonable that if present state-of-the-art digital computer techniques can synthesize speech, it should be possible to determine the intelligibility of speech using computer processing. Hopefully the approach used in this thesis will provide at least a basis for development of a computerized method for measuring intelligibility that will prove to be sufficiently accurate and simple to replace the human listener method. The ultimate development of this type technique would provide a device which could be hooked to the communications system under test and have a meter which would indicate the intelligibility of the system on a real-time basis.

I am indebted to Major Joe Carl, my advisor, for his guidance, suggestions, advice, and encouragement during the preparation of this thesis. I would like to express my appreciation to Captain Mazzie and Mr. William Hall, Jr. of the Analog/Hybrid Systems Branch of the ASD Computer Center for the many hours they spent on the preliminary processing of the analog speech data. My thanks to Dr. Oestreicher and Richard McKinley of the Aerospace Medical Research Laboratory for allowing me to use their anechoic chamber to prepare the voice data tapes and their computer terminal to develop the computer algorithms and process

the data. I would also like to express my appreciation to Captain John Bauer for providing all the subjective listener data used as a basis for comparison in this thesis.

Wayne R. Beeson

Table of Contents

	<u>Page</u>
Preface	ii
List of Figures	v
List of Tables	vi
Abstract	vii
I. Introduction	1
Background	1
Approach	5
Objective	7
Scope	8
II. Data Acquisition	9
Acquisition Procedure	9
III. Analog to Digital Conversion	12
Frequency Analysis	14
IV. Digital Signal Processing	16
Data Compression	16
Gray Scale Spectrogram	18
Word Location Technique	18
Master Tape Processing	20
V. Cross-Correlation and Mean Squared Error Calculation	23
VI. Results	26
VII. Conclusions and Recommendations	32
Bibliography	34
Appendix A: Sequence Chart for Intelligibility Prediction	36
Appendix B: Computer Programs	41
Vita	50

List of Figures

<u>Figure</u>		<u>Page</u>
1	Process Used to Calculate Articulation Index (AI)	4
2	Timing Tone Generator	13
3	Word Spectrogram	19
4	Digital Speech Spectrogram and Associated Matrix Array.	22
5	Matrix of Cross-Correlation Values for a Word at S/J 7	27
6	Plot of Mean Squared Error Versus S/J	28
7	Plot of Percent of Listener Wrong Answers Versus S/J Ratio	29
8	Scatter Plot of Mean Squared Error Versus Percentage of Listener Wrong Answers	31
9	Sequence Chart for Master Tape Processing Program (Plate 1)	37
10	Sequence Chart for Master Tape Processing Program (Plate 2)	38
11	Sequence Chart for Cross-Correlation and MSE Program (Plate 1)	39
12	Sequence Chart for Cross-Correlation and MSE Program (Plate 2)	40

List of Tables

Table		<u>Page</u>
I	Diagnostic Rhyme Test	10
II	Speech Frequencies	17
III	Overprint Symbols for Speech Spectrograms	18

Abstract

A method of predicting speech intelligibility using computer algorithms is presented. Diagnostic Rhyme test number four was used to measure speech intelligibility using a subjective listener test and these results were used as a basis for comparison with the intelligibility predictions made by the computer algorithm. An audio recording of a speaker reading the Diagnostic Rhyme test was made. This recording was run through a General Electric radio system and varying amounts of noise were added. The output of the radio system was recorded, providing a copy of the input word corrupted by both additive noise and radio system distortion effects. Both the input recording and the noisy output recording were digitized by sampling the analog waveforms at a 10 kilohertz rate. These digital samples were converted to a frequency format by windowing the time samples with a rectangular window 128 time samples in length and processing them using Fast Fourier transform techniques. This procedure simulated running the analog speech signal through a bank of contiguous narrow bandpass filters covering the range of 0 to 5 KHz, with center frequencies 78 Hz apart. The output of this process was a matrix array, corresponding to each word from the tape, of amplitude values 200 time windows long and divided into 64 frequency bands. These 64 frequency bands were then combined into 1/3 octave groups to model the frequency sensitivity of the average human ear, which reduced the matrix array to 16 frequency bands. This processing of the analog signal was used to model the preprocessing which occurs in the human ear. A comparison between each word from the input tape and the noisy output tape was then made using a weighted mean squared error

calculation. This comparison was conjectured to provide a difference measure which is inversely related to intelligibility. This comparison was used to represent how intelligible the input received from the inner ear is to the brain.

Comparison of the intelligibility results from the human listener tests with the computer processing method outlined above gave a Pearson's Correlation Coefficient value of 0.74 which indicates the computer prediction accounted for 55% of the variance in the listener error scores.

AN ALGORITHM FOR DETERMINING
SPEECH INTELLIGIBILITY

I. Introduction

This work is in response to a need identified by Air Force Communications Service (AFCS). There have been numerous studies in the area of machine prediction of speech intelligibility; however, the Air Force planners who requested this work are still using human listeners to determine intelligibility. They either think that available computer methods do not produce sufficiently accurate results or that the computer schemes are too complex and difficult to apply to their specific problems. The intent of this work is to take applicable techniques from work that has already been done and combine them to develop a simplified, accurate method to evaluate voice intelligibility with a computer.

Background

The oldest method for determining the intelligibility of speech is a subjective method that involves trained speakers and listener panels that directly score the percentage of speech that is intelligible. This method is still considered the most reliable way to measure intelligibility because it produces repeatable results. The disadvantages of the subjective method are the considerable cost, large number of manhours, and specialized facilities and equipment required.

An early attempt to simplify the procedure for determining the intelligibility of speech involved calculation of the mean squared error (MSE) between an audio waveform and the same audio signal corrupted by noise. This process uses the procedure given by Equation 1.

$$\text{MSE} = \frac{1}{2T} \int_{-T}^T |x(t) - y(t)|^2 dt \quad (1)$$

Such an approach does not yield acceptable estimates of speech intelligibility. This failure is attributed to the fact that vowels contain more power than consonants in speech, but consonants are more important in determining intelligibility than vowels (Ref 15:277). This method of intelligibility measurement is no longer in use due to these shortcomings.

One of the currently popular methods of automating intelligibility prediction is use of the Articulation Index (AI). One method of calculating the AI is by transforming the speech signal into an electrical signal and then passing it through a set of contiguous bandpass filters each 1/3 octave wide. The voltage output of each of these filters is used to calculate a root mean square (RMS) voltage as shown by Equation 2

$$\text{RMS} = \frac{1}{2T} \int_{-T}^T x^2(t) dt \quad (2)$$

The noise that is affecting the system is passed through this same set of filters and a root mean square noise voltage in each filter bandpass is calculated. The value of the noise RMS voltage is subtracted from the speech RMS voltage for each filter. If this difference is 30 or more decibels, it is assigned a value of 30. If the difference falls in the range of 0 to 30 decibels, the actual decibel value is assigned. If the difference is 0 or a negative value, it is assigned a value of 0. These values for each filter are then multiplied by weighting factors for each of the different frequency bands. These products are then

added together and their sum is the AI (Ref 1:6-15). This process is illustrated by the block diagram in Figure 1.

The AI can be calculated by programming the previous procedure into a computer routine. These programs are in common use today and are an acceptable predictor of intelligibility as long as the noise present is additive white Gaussian noise. Colored noise and multiplicative noise require a complete recalibration of the AI system to give good results (Ref 7:2). This illustrates that the type of noise present must be known exactly and corrected for to maintain the accuracy of this index. When this method of intelligibility prediction is applied to a digital voice system or a system with quantization noise present, there is no correction of recalibration that will provide acceptable intelligibility estimates (Ref 7:2-3).

A recent development in the field of automated voice intelligibility prediction is the use of linear predictive coding (LPC). LPC derives its name from the predictive process it is based on which states: given P samples of a speech signal, the next sample can be predicted approximately by a linear function of the P known samples (Ref 6: A-3). LPC models the vocal tract as an all-pole digital filter and estimates the filter parameters (predictor coefficients) using the time domain speech waveform. This model of the voice tract assumes the vocal tract model to be a time-varying filter with parameters changing slowly enough so they can be considered fixed over a specified time interval. It accounts for the glottal volume flow and radiation of sound from the mouth in addition to vocal tract sounds (Ref 7:3-4). The most popular way to estimate linear prediction coefficients (a_i) is the autocorrelation method. This method involves time sampling an analog speech signal and

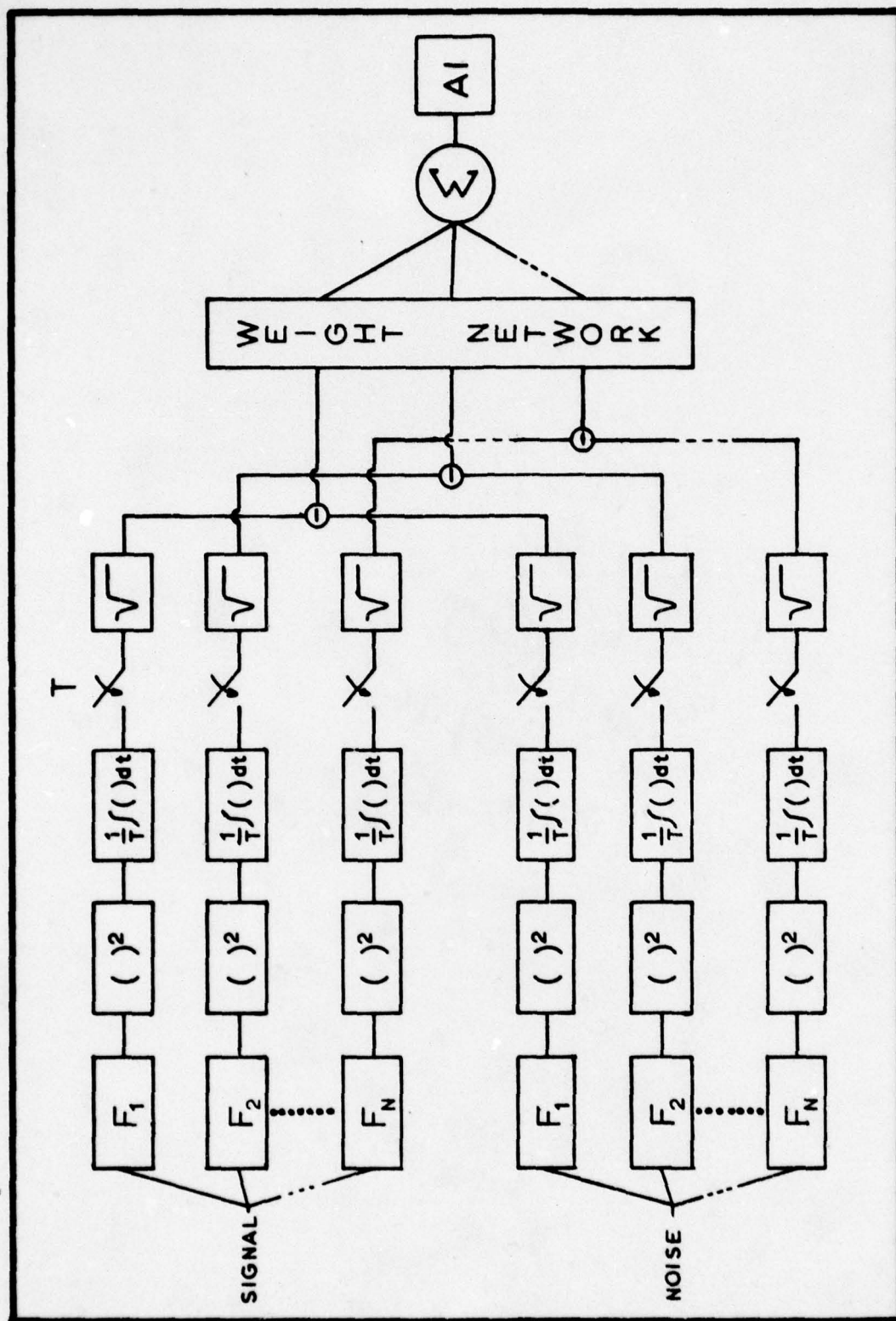


Figure 1. Process Used to Calculate Articulation Index (AI)

windowing these time samples, usually with a Hamming window 256 time samples long. These windowed speech samples are used to calculate a linear prediction of residual energy both for an undistorted speech signal and for the same speech signal after it has been corrupted by additive noise. These residual energy terms are then compared with the actual energy terms and a distance measure is derived from these comparisons. The calculation of these residual energy terms and their comparison is a long and involved mathematical process presented in detail by Hartmann (Ref 6:24-33).

The use of LPC techniques overcomes the disadvantages of sensitivity to the type of noise present that affects AI. LPC intelligibility predictions give a good correlation with listener scores when averaged over 50 or more words (Ref 6:18-19). The disadvantage of LPC is that it involves a large number of computer computations and consequently consumes a great deal of computer time to analyze a small amount of speech. A second disadvantage of this method is that it requires very close synchronization of the words on the undistorted tape with the same words on the tape containing additive noise. This requirement for exacting synchronization makes it necessary to employ very specialized taping equipment to make this process work (Ref 6: 18,20).

Approach

The approach to computer evaluation of speech intelligibility used in this thesis combines some features of the Articulation Index calculation, the linear predictive coding method, and the mean squared error calculation.

The human auditory system performs multiple stages of preprocessing on an audio signal before it reaches the brain. Therefore, it seems

reasonable to assume that if the processes occurring in the ear and brain can be modeled, it will be possible to make the same type intelligibility determination as the human. The first step in doing this is to model the preprocessing which occurs in the ear.

To model the action of the ear drum in converting sound pressure variations to vibration and the middle ear which transmits these as a varying mechanical vibration to the inner ear, a tape recorder was used. The recorder converts sound pressure variations into an appropriate, continuously varying analog signal.

In the inner ear (cochlea) the mechanical vibration variations undergo the next stage of processing. This process is quite complex, but it appears to involve excitation of the neurons at the base of the hair cells inside the cochlea due to movement of the hair cells. This movement is a result of the mechanical vibration coming from the middle ear causing the fluid in the cochlea to move the hair cells. Since the cochlea is apparently a frequency analyzing device, a model for the inner ear should present the signal in a frequency format (Ref 10). The model for the cochlea used in this thesis consists of sampling the analog waveform from the tape at the Nyquist rate and running these samples through a bank of contiguous bandpass filters. The output of these filters are grouped into $1/3$ octave bands to simulate the sensitivity of the ear. This changes the analog waveform into a frequency format.

Kabrisky proposed that the cortex of the brain is capable of performing a two-dimensional cross-correlation of a test image with a stored pattern (Ref 9: 47-57). This theory about the visual system was extended to the auditory system by Dailey and Sutton (Ref 3). Assuming

the brain performs this two-dimensional cross-correlation between the input signal from the cochlea and phonemes stored in memory and picks the largest correlation value to indicate what phoneme was heard, this process must be modeled. Since the undistorted phoneme stored in memory would have to be in the same format as the incoming signal from the cochlea, a possible model of this process would be to compare the undistorted input word, preprocessed by the models of the ear, with the same word imbedded in noise and run through the same processing. The correlation process will only determine a measure of the difference between a word and its corrupted form, so it appears that a mean squared error calculation can simulate this correlation satisfactorily. This mean squared error will be weighted because of the grouping of filter outputs occurring in the cochlea model.

Objective

The object of this research is to explore the possibility of developing a computer program that will give a reasonably accurate prediction of the intelligibility of speech. This system, if successful, will be used by people with varying degrees of computer support available to them. For this reason the main idea was to keep the procedure simple and automate it as much as possible. Another consideration was to minimize the computer memory and central processor time required for the processing so people can get the program through a busy, time shared computer in a reasonable amount of time. The final goal was to eliminate the need for any elaborate or unique equipment to make or process the audio tape.

Scope

The scope of this project is limited to developing a computer program that will model the actions of the ear on sound waves and apply one possible comparison scheme to model the action of the brain on the processed audio signal. Section II outlines the procedures used to make audio tapes that are used for both human listener and computer intelligibility testing. Section III details the initial computer processing of these audio tapes to sample them at the Nyquist rate and perform a Fast Fourier Transform (FFT) on these time samples. Section IV describes how the original matrix array of amplitude values, produced by the FFT process, was compressed so it would closely approximate the way the ear processes sound data. A method of representing each word by a speech spectrogram and using this to locate the word exactly in a group of time samples of the input wave form is discussed. Section V deals with the cross-correlation method used to locate a word which is imbedded in noise. It evaluates how much the word has been distorted by the additive noise using a weighted mean squared error comparison between the word before the noise is added and the same word plus additive noise. The last two sections show the results of this procedure and make some recommendations for further work in this area.

II. Data Acquisition

The data used in these tests for intelligibility was the Diagnostic Rhyme Test Number IV (DRT-IV). DRT-IV is composed of 58 rhyming word pairs with each word pair designed to test for one of six speech attributes. There are eight rhyming pairs in the list which check for each attribute. The six attributes tested for are voicing, nasality, sustention, sibilation, graveness, and compactness (Ref 15:15-21). The words that test for these attributes are separated by ten pairs of filler words. The DRT-IV used in these tests is shown in Table I and the words which test for each of the speech attributes are identified.

Acquisition Procedure

The data acquisition consisted of a male speaker reading one word of each rhyming word pair from DRT-IV and recording these words on one track of a stereo tape recorder. The other track of the stereo tape was used to record one kilohertz tones which are used for timing references in subsequent processing of the audio tape. The recorder used was a reel-to-reel Sony Model 850 which gave a reasonably high quality of audio reproduction. Recording of the DRT-IV words on tape was used to model the action of the outer ear which converts the pressure variations of sound into an analog signal format.

In recording the test audio tapes, two different male speakers were used to reduce the possible effect of a speaker's regional accent affecting the intelligibility results. The first speaker had a southern accent (Arkansas) and the second had very little regional accent (Minnesota). Four master tapes were made of DRT-IV, two by the first speaker and two by the second speaker. These four master tapes were

Table I
Diagnostic Rhyme Test

DRT IV-(2)

PEST - TEST	-(filler)-	FAN - PAN
VAULT - FAULT	-(voicing)-	CHOCK - JOCK
DUES - NEWS	-(nasality)-	NOTE - DOTE
VEE - BEE	-(sustention)-	TICK - THICK
THANK - SANK	-(sibilant)-	CARE - CHAIR
ROD - WAD	-(graveness)-	DONG - BONG
SO - SHOW	-(compactness)-	YOU - RUE
LID - RID	-(filler)-	REEK - LEAK
DENSE - TENSE	-(voicing)-	GAFF - CALF
BOSS - MOSS	-(nasality)-	BOMB - MOM
FOO - POOH	-(sustention)-	DOUGH - THOUGH
ZEE - THEE	-(sibilant)-	GILT - JILT
FAD - THAD	-(graveness)-	PENT - TENT
HOP - FOP	-(compactness)-	YAWL - WALL
ROW - LOW	-(filler)-	LOOT - ROOT
GIN - CHIN	-(voicing)-	VEAL - FEEL
BEND - MEND	-(nasality)-	NAB - DAB
CHAW - SHAW	-(sustention)-	BON - VON
JUICE - GOOSE	-(sibilant)-	SOLE - THOLE
PEAK - TEAK	-(graveness)-	THIN - FIN
BAT - GAT	-(compactness)-	KEG - PEG
ROCK - LOCK	-(filler)-	LONG - WRONG
GOAT - COAT	-(voicing)-	TUNE - DUNE
MIT - BIT	-(nasality)-	MEAT - BEAT
THEN - DEN	-(sustention)-	SHAD - CHAD
GAUZE - JAWS	-(sibilant)-	GOT - JOT
NOON - MOON	-(graveness)-	DOLE - BOWL
KEY - TEA	-(compactness)-	DILL - GILL
RAMP - LAMP	-(filler)-	LEND - REND

played into the input of a General Electric (GE) preliminary development model, spread-spectrum radio transmitter. The modulated radio signal was transmitted to a hybrid summer where it was mixed with additive noise from a pseudo-random, matched spectrum noise generator. The output of the hybrid summer was then fed into the companion receiver of the GE transmitter and the output audio tapes were recorded at the audio output stage of the receiver. The noise generator was used to simulate intentional jamming of the radio link. Output tapes were made at eight different signal to jammer (S/J) levels. The four input tapes were each used twice as the input when making the different S/J level output tapes. The eighth output tape had the lowest S/J ratio and the signal was too low to be usable for testing, so this tape was discarded. The remaining tapes have S/J levels numbered from one through seven. The highest S/J level is number one and the S/J ratio decreases as the number increases with tape number seven representing the lowest S/J ratio. The actual S/J levels associated with these numbers are classified Secret. If the actual S/J levels are desired, this information is given in the classified portion of Captain Bauer's thesis (Ref 2).

All recordings of the output of this system were made on new Scotch, 1/4 inch tape using a Sony Model 850 recorder. These tapes were recorded at a tape speed of 7½ inches per second.

The GE communications system used in this test was being evaluated by a fellow student for intelligibility using human listener tests (Ref 2). This provided a convenient way to obtain human listener intelligibility data to compare to the intelligibility predictions made by the computer method presented in this thesis.

III. Analog to Digital Conversion

The initial processing of the four master DRT-IV input tapes and the seven output tapes with different S/J additive noise ratios, was done by the Analog/Hybrid Systems Branch of the Aeronautical Systems Division (ASD) Computer Center.

When the audio tapes were made, the words from DRT-IV were recorded on the right channel of a stereo tape recorder and a one kilohertz (KHz) tone, 1/2 second long, was recorded on the left channel. These one KHz tones were machine generated with the apparatus shown in Figure 2. Every time the tone sounded, the next word from the DRT-IV list was recorded within 2½ seconds after the tone. The tones were spaced seven seconds apart so there was at least 4½ seconds after each word before the next tone. The Analog/Hybrid Branch played each tape back and low pass filtered it to 2.5 KHz and fed this into the Comcor Ci-5000/6 analog computer. The computer sampled the input at the Nyquist rate of 5 KHz. Using 2.5 KHz as the upper cutoff frequency was necessary because the bandwidth of the amplifiers in the analog computer was limited to this value. Since it was desired to analyze the speech input over a range of zero to 5 KHz, it was necessary to analyze the tape by playing it back at a tape speed of 3 3/4 inches per second, half the recording speed, to give the effect of low pass filtering to 5 KHz and sampling at a 10 KHz rate at the original recording speed. This makes it possible to evaluate the speech signal over the desired frequency range in spite of the limitations imposed by the computer's amplifier bandwidth.

The input speech signal was amplified to approximately 100 volts prior to processing to provide a sufficient voltage swing to utilize

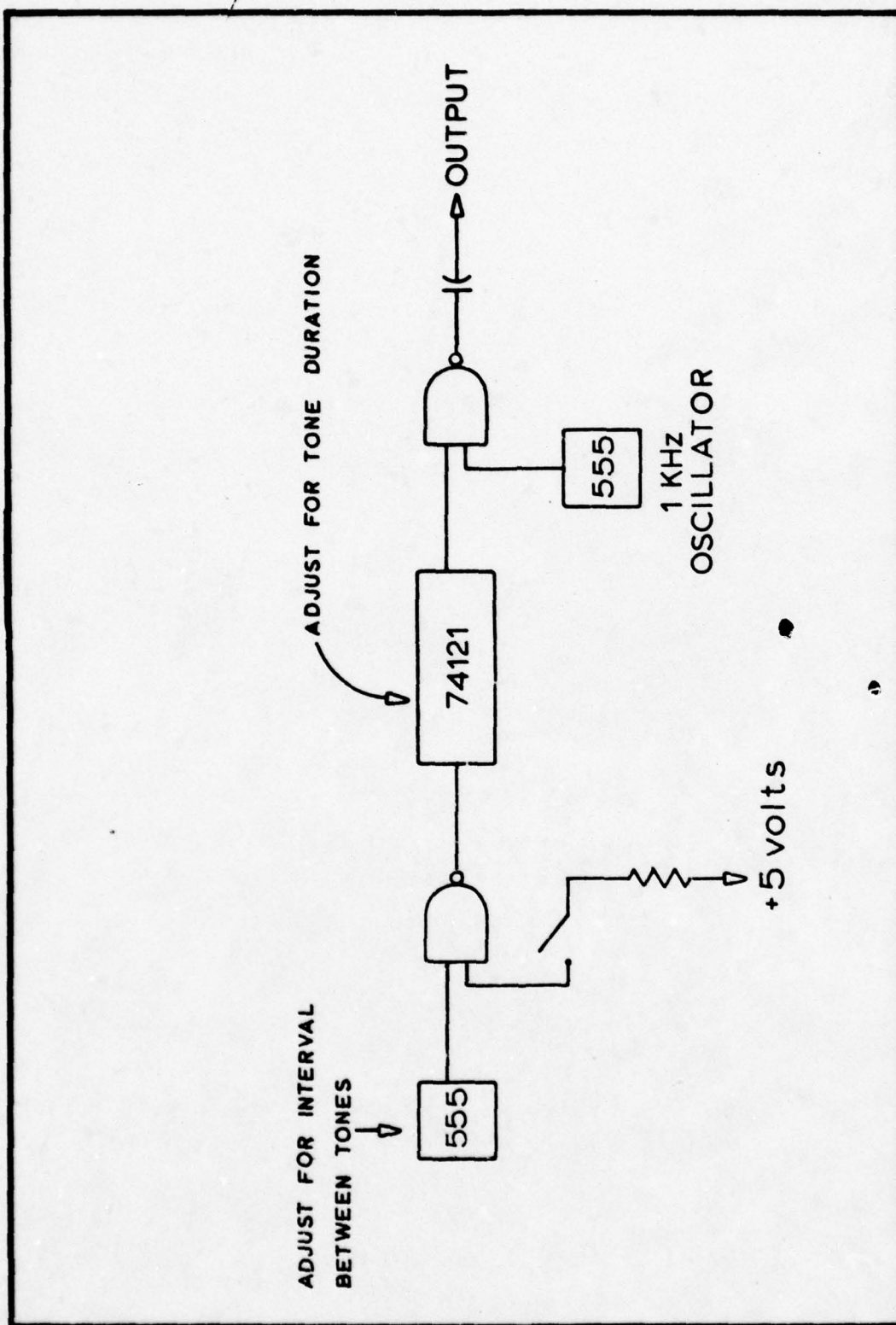


Figure 2. Timing Tone Generator

the accuracy possible with the 11 bit analog to digital converters in the computer. These 11 bit numbers are a binary representation of a 4 digit decimal number. These numbers give the voltage level, between 0 and 100 volts, of the analog waveform each time the waveform is sampled by the digital computer. The 1 KHz tones recorded on the left channel of the audio tapes were used to trigger the sampling equipment in the computer. When the tone occurs, the computer starts sampling the input audio waveform and continues sampling for $2\frac{1}{2}$ seconds, then stops until the next tone occurs. The word from DRT-IV is contained somewhere in this $2\frac{1}{2}$ second sampling interval and will be located exactly by subsequent processing.

Frequency Analysis

The proposed model for the inner ear requires that the digitized analog speech data be represented in the equivalent frequency domain. Fast Fourier Transforms (FFT) techniques were used to convert the digitized data into a frequency representation format (Ref 5:41-52). The actual data conversion involves grouping the digitized time samples into groups of equal length (windowing) and applying FFT techniques to these window groupings to simulate a bank of narrow bandpass filters. The size and shape of the window is based on the desire to have a wide band analysis while retaining reasonable time resolution. The methods for doing this are discussed in detail by Neyman (Ref 12:17-18). The window used is rectangular and 128 time samples long. The digitized data was processed using this window size by an Analog/Hybrid Branch program called AMPSPC. This program gave 64 discrete amplitude values, each corresponding to a 78.125 Hz frequency segment located in the range of 0 to 5 KHz and covering a time window of 12.8 milliseconds

(128 time samples at 10 KHz sample rate). This produced a 64 x 200 matrix array of amplitude values. Each of these matrix arrays contains one word from the DRT-IV audio tape. The matrix arrays were written on a nine track ASD Computer Center library tape (L-tape) and stored for later processing by the CDC-6600 computer.

IV. Digital Signal Processing

Each word of the DRT-IV is now contained in a 64 x 200 matrix array stored on a computer L-tape. Each element of this array represents the signal amplitude to four decimal place accuracy. This signal representation corresponds to an analog speech signal that has been run through a bank of 64 bandpass filters with center frequencies 78.125 Hz apart and an upper cutoff of 5 KHz.

Data Compression

In general the human ear is not a linear receiving device. In order to model the nonlinear frequency response of the ear it was necessary to restructure the digital data so it would approximate the ear's unusual sensitivity to frequency change. The six lowest frequency bands of the matrix array were left unchanged. This group has center frequencies of 78.125, 156.250, 234.375, 312.500, 390.625, and 468.750 Hz. All higher frequencies, up to 5 KHz, are grouped into approximately 1/3 octave ranges and the energy content of each group is the sum of the individual array elements contained within that group. This restructuring produced 16 frequency dependent amplitudes from the original 64. The frequency groupings that produce these 16 values are shown in Table II.

Adding the energy of each array element within a 1/3 octave group compensated for the lower amplitude of sound harmonics produced by the vocal chords at high frequencies. This eliminated the need to use the standard preemphasis technique of increasing the signal magnitude by six decibels per octave above 350 Hz (Ref 14:311).

Table II

Speech Frequencies

Center Frequency Original Data	Center Frequency Original Data	Center Frequency Original Data
78.125	78.125	2578.125
156.250	156.250	2656.250
234.375	234.375	2734.375
312.500	312.500	2812.500
390.625	390.625	2890.625
468.750	468.750	2968.750
546.875	585.940	3046.875
625.000		3125.000
703.125	742.188	3203.125
781.250		3281.250
859.375	898.440	3359.375
937.500		3437.500
1015.625	1132.810	3515.625
1093.750		3593.750
1171.875		3671.875
1250.000		3750.000
1328.125		3828.375
1406.250	1445.310	3906.250
1484.375		3984.375
1562.500		4062.500
1640.625		4140.625
1718.750		4218.750
1796.875	1793.380	4296.875
1875.000		4375.000
1953.125		4453.125
2031.250		4531.250
2109.375		4609.375
2187.500	2226.560	4687.500
2265.625		4765.625
2343.750		4843.750
2421.875		4921.875
2500.000		5000.000

2812.500

3554.690

4453.125

Gray Scale Spectrogram

It is desirable to be able to see a spectrogram of the word when working with the 16 x 200 matrix array that results from the previous compression procedure. This aids in quickly locating where the word occurs in the 200 time windows and determining the length of the word. A convenient method for creating a gray scale spectrogram using computer overprint symbols was developed by Neyman and is used here (Ref 12:22-24). Table III shows the overprint symbols used to create the spectrogram and Figure 3 shows what the actual spectrogram of a word looks like.

Table III

Overprint Symbols for Speech Spectrograms

Number of Overprints	LEVEL OF DARKNESS									
	0	1	2	3	4	5	6	7	8	9
1			+	x	x	x	x	x	x	x
2					-	+	0	0	0	0
3								-	-	#
4									+	+
5										*

Word Location Technique

In order to use the matrix array of amplitude values for comparison purposes it is necessary to locate the word in the 200 time windows and save only the part of the array containing the word. This location process was accomplished by thresholding each value in the matrix array at 1.5 and saving the part of the array where the values exceeded this level. A filtering process was included in the program to eliminate a

	+	+++	36
X	++	+++	37
X	+++	+++	38
X	+++	+++	39
X	+++	+++	40
X	+++	+++	41
X	+++	+++	42
X	+++	+++	43
X	+++	+++	44
X	+++	+++	45
X	+++	+++	46
X	+++	+++	47
X	+++	+++	48
X	+++	+++	49
X	+++	+++	50
X	+++	+++	51
X	+++	+++	52
X	+++	+++	53
X	+++	+++	54
X	+++	+++	55
X	+++	+++	56
X	+++	+++	57
X	+++	+++	58
X	+++	+++	59
X	+++	+++	60
X	+++	+++	61
X	+++	+++	62
X	+++	+++	63
X	+++	+++	64
X	+++	+++	65
X	+++	+++	66
X	+++	+++	67
X	+++	+++	68
X	+++	+++	69
X	+++	+++	70
X	+++	+++	71
X	+++	+++	72
X	+++	+++	73
X	+++	+++	74
X	+++	+++	75
X	+++	+++	76
X	+++	+++	77
X	+++	+++	78
X	+++	+++	79
X	+++	+++	80
X	+++	+++	81
X	+++	+++	82
X	+++	+++	83
X	+++	+++	84
X	+++	+++	85
X	+++	+++	86
X	+++	+++	87
X	+++	+++	88
X	+++	+++	89
X	+++	+++	90
X	+++	+++	91
X	+++	+++	92
X	+++	+++	93
X	+++	+++	94
X	+++	+++	95
X	+++	+++	96
X	+++	+++	97
X	+++	+++	98
X	+++	+++	99
X	+++	+++	100

Figure 3. Word Spectrogram

noise spike, occurring outside the word, from being mistaken for part of the word.

Master Tape Processing

The plan was to sample the analog speech data, convert it to a frequency format representation, compress the resulting 64 frequency divisions to 16, locate the word exactly in the 200 time windows, and store only the portion of the 200 column array where the word occurs for use in subsequent processing. These steps were written into a computer program for use in the CDC-6600 computer.

This program was used to process the four master DRT-IV tapes that were used as inputs to the GE radio system. Each master tape was run through this program and the column number where each word started, the number of columns (M) occupied by the word, and the 16 x M array of amplitude values containing the word were recorded on a second L-tape.

An alternate method for locating the word in the 200 time windows was tried, to establish a basis for comparing the effectiveness of the previous procedure. The matrix array was first normalized using Equation 3

$$\hat{a}_{i,j} = \frac{a_{i,j}}{\left[\sum_{i=1}^L \sum_{j=1}^{16} a_{i,j}^2 \right]^{1/2}} \quad (3)$$

where $\hat{a}_{i,j}$ = normalized array element. The normalized array was then thresholded at an appropriate level and the computer program predicted where the word was located at in the 200 column matrix. It was easy to see where the word occurred within the 200 column matrix by looking

at the accompanying spectrogram, Figure 4, so this was used as a means of evaluating the two computer techniques given above for locating the word. This comparison showed that by normalizing the data in the array prior to doing a computer search for the word, the computer program frequently failed to correctly locate the word. When the computer search for the word was done without normalizing the matrix array, it found the word accurately every time. No explanation can be offered to explain why normalizing the array data caused the word location program to fail.

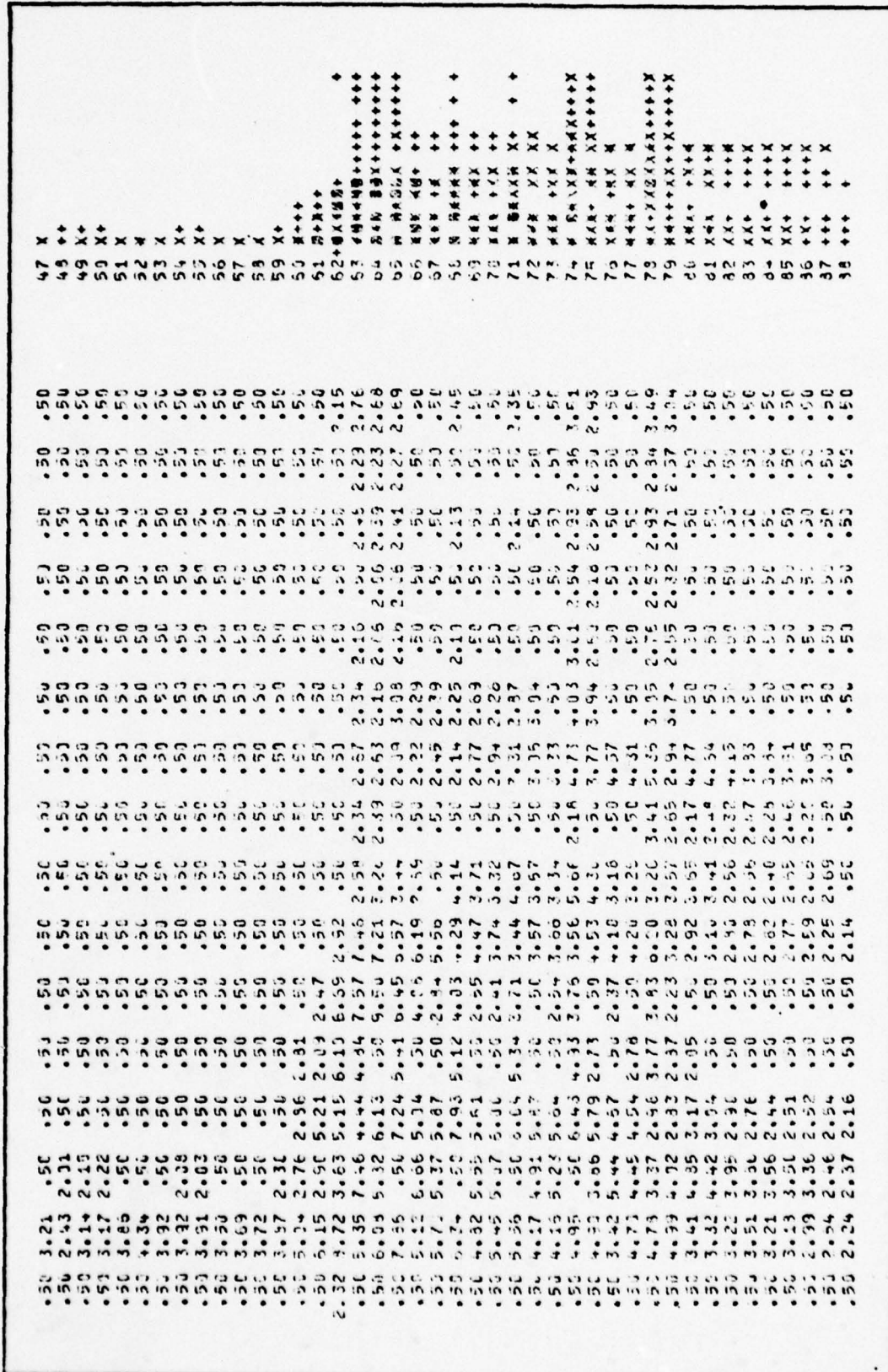


Figure 4. Digital Speech Spectrogram and Associated Matrix Array

V. Cross-Correlation and Mean Squared Error Calculation

The seven noisy tapes made at the audio output of the GE radio under test must now be compared with the master input tape from which each was made. The input tape is compared with the output tape and the mean squared error between each input word and its output plus additive noise is calculated.

To locate where the word occurred on the noisy output tape, it was necessary to perform a cross-correlation between the input word array and the 16 x 200 output array containing the same word imbedded in noise. The length of the word and the part of the array containing the word from each master tape were previously recorded on a computer L-tape. This L-tape is read a word at a time and cross-correlated with the corresponding 16 x 200 array on the L-tape containing the noisy words. The length of the word is read first and that number subtracted from 200 to find the number of cross-correlations that must be performed. Next the arrays are read into the computer core memory and a cross-correlation is performed with the first column of the word from the master tape lined up with the first column of the 16 x 200 noisy array. After each cross-correlation value is determined, the array containing the word from the master list is shifted one column to the right with respect to the noisy array. When all the cross-correlations have been performed for that word, the largest value computed indicates the point where the input word and the noisy output word were aligned. The equation used to compute each of these cross-correlations (P) is

$$P(\tau) = \sum_{i=1}^L \sum_{j=1}^{16} A_{i,j} B_{i,j} \quad (4)$$

where $L = 200$ - length of word read from master tape

$A_{i,j}$ = element of word array from master tape

$B_{i,j}$ = element of word array from noisy tape

When the maximum cross-correlation value has located the word (that is, once the value of τ_0 is known such that $P(\tau_0)$ is a maximum) in the 16×200 array containing the additive noise, there is enough information available to calculate the mean squared error (MSE) between the two words using the formula

$$MSE = \sum_{i=1}^L \sum_{j=1}^{16} (A_{i,j})^2 - 2 P(\tau_0) + \sum_{i=\tau_0}^{L+\tau_0} \sum_{j=1}^{16} (B_{i,j})^2 \quad (5)$$

To compute the MSE the maximum cross-correlation value determined in the previous operation is multiplied by two and is the middle term of the MSE equation. Since the exact location of the word in the 16×200 noisy array is now known, the last term of the MSE equation can be calculated using this information to square the elements of the array containing the word and sum these squares. The L-tape containing the array of the word from the master input list is used to calculate the first term in the MSE equation by squaring each element of the array and then summing the squares.

The average mean squared error for the DRT-IV list recorded at each of the seven different S/J ratios was determined by summing the MSE for all 58 words in the list and dividing this sum by 58. This gives an average MSE corresponding to each S/J level to be used in deriving an estimate of the intelligibility at that S/J level.

The brain "expects" to see uncorrupted phoneme groups (words) characterized in a certain way. Assuming that the models used here give a

reasonably accurate representation of the preprocessing which occurs in the ear, it should now be possible to measure the difference between what the brain expects to hear and what it actually hears. It is conjectured that this difference is inversely related to the intelligibility of what is heard. To model this process the MSE calculation provides a measure for determining the difference between a word and that same word after it has been corrupted by noise. It is assumed this distance measure can now be related to intelligibility.

VI. Results

The first computer program discussed in section IV was designed to locate each word in the 200 time windows. This program worked perfectly as long as the tapes containing the words had a low noise level compared to the signal amplitude of the word. It was also necessary to amplify the peak levels within each word to at least 75 volts prior to digital sampling.

The second program, discussed in section V, was designed to first locate the word in the 200 time windows on the tape with additive noise by a cross-correlation with the same word without noise. This cross-correlation provided a sharp peak with an amplitude well above the other values to indicate when the two words were aligned. This distinct peak occurred even at the lowest S/J ratio. This can be seen by looking at the cross-correlation values for a word at the lowest S/J level, Figure 5.

The second part of the program described in section V is used to calculate the mean squared error between the input word and the same word after it passes through the radio system and is corrupted by noise. Figure 6 shows a plot of mean squared error values versus the seven different S/J ratios. The increase in MSE is approximately linear as the S/J ratio decreases.

Figure 7 shows the average number of errors made by the 10 people who listened to the noisy tapes at each S/J ratio.

A scatter plot of human listener error scores versus mean squared error values is shown in Figure 8. This plot displays the data used to calculate Pearson's Correlation Coefficient.

RECORD NUMBER 7 HAS A MAXIMUM CROSSCORRELATION VALUE OF 3443.45

VALUES OF EACH CROSSCORRELATION

2283.15	2286.36	2307.44	2350.55	2375.35
2439.34	2427.66	2403.39	2442.67	2452.12
2512.39	2466.74	2523.84	2494.79	2477.39
2507.18	2549.63	2541.09	2499.52	2529.81
2528.40	2584.06	2523.57	2485.37	2472.34
2529.14	2622.18	2585.15	2574.93	2561.39
2574.73	2585.51	2636.59	2636.58	2602.56
2639.36	2595.76	2535.39	2552.40	2561.04
2587.32	2620.85	2637.04	2639.65	2648.63
2539.60	2563.31	2521.55	2575.39	2529.16
2583.34	2538.89	2536.94	2525.86	2562.09
2491.09	2501.41	2523.95	2483.81	2403.49
2371.19	2310.93	2310.50	2283.24	2232.59
2231.78	2204.32	2222.35	2197.52	2133.11
2133.78	2117.15	2074.58	2060.74	2067.75
2056.84	2053.27	2074.08	2103.98	2077.40
2105.43	2109.71	2136.01	2225.55	2279.52
2321.64	2415.90	2449.44	2520.54	2561.87
2621.29	2652.15	2687.30	2774.00	2886.51
3106.15	3095.53	3145.38	3213.07	3246.50
3314.19	3363.76	3386.61	3397.87	<u>3443.45</u>
3404.21	3322.75	3301.64	3297.31	3224.50
3195.33	3116.08	3029.27	2975.67	2937.41
2847.04	2764.71	2714.32	2674.48	2651.47
2724.31	2729.74	2687.19	2682.50	2622.32
2542.56	2644.09	2633.69	2573.14	2548.18
2508.40	2548.13	2552.30	2513.10	2499.71
2558.47	2573.17	2595.04	2571.22	2558.32
2621.04	2605.42	2636.91	2606.57	2612.94
2613.22	2633.35	2644.85	2699.65	2657.45
2651.23	2627.83	2619.78	2652.16	2625.06
2611.55	2633.13	2595.43	2605.58	2624.91

Figure 5. Matrix of Cross-Correlation Values for a Word at S/J 7

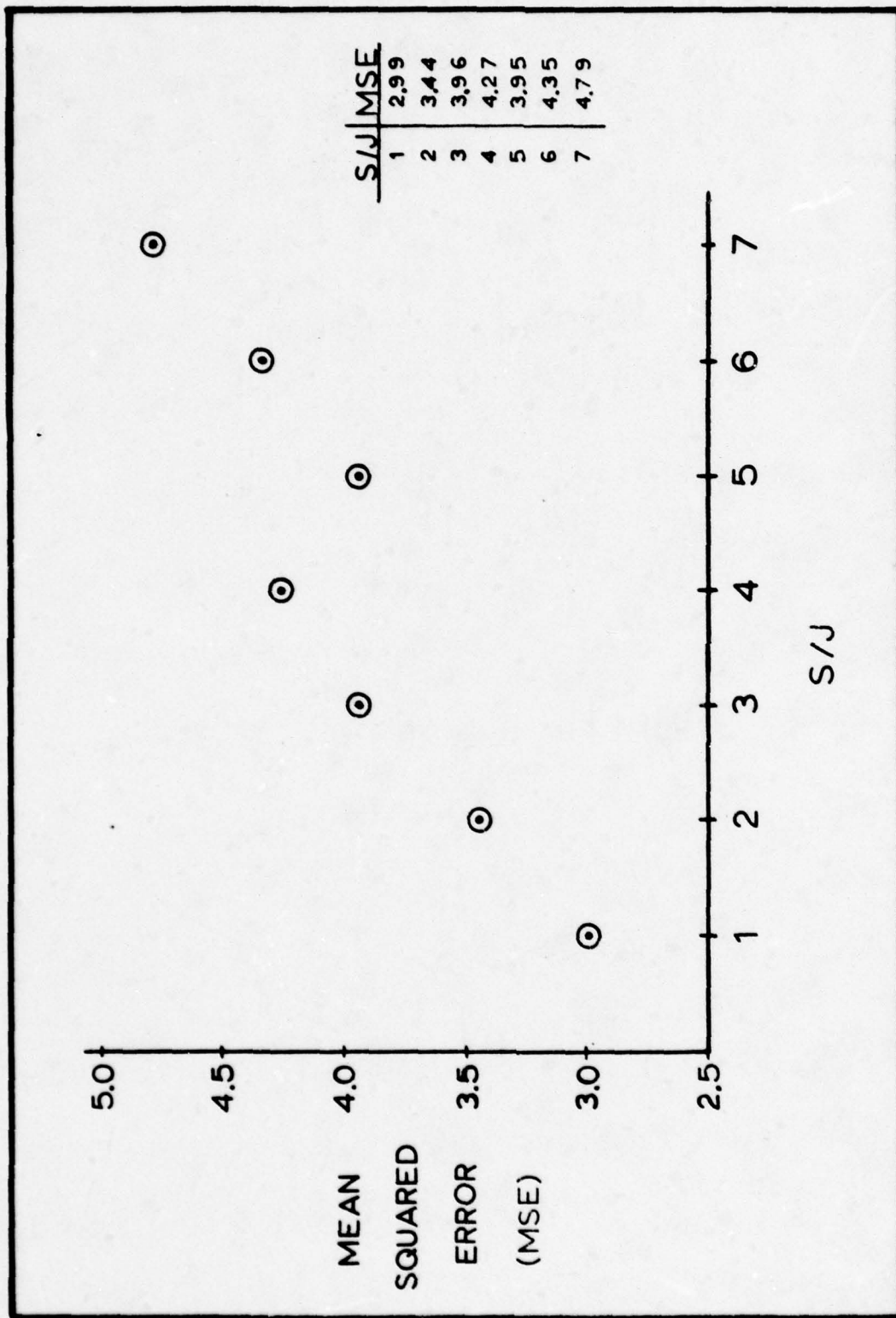


Figure 6. Plot of Mean Squared Error Versus S/J

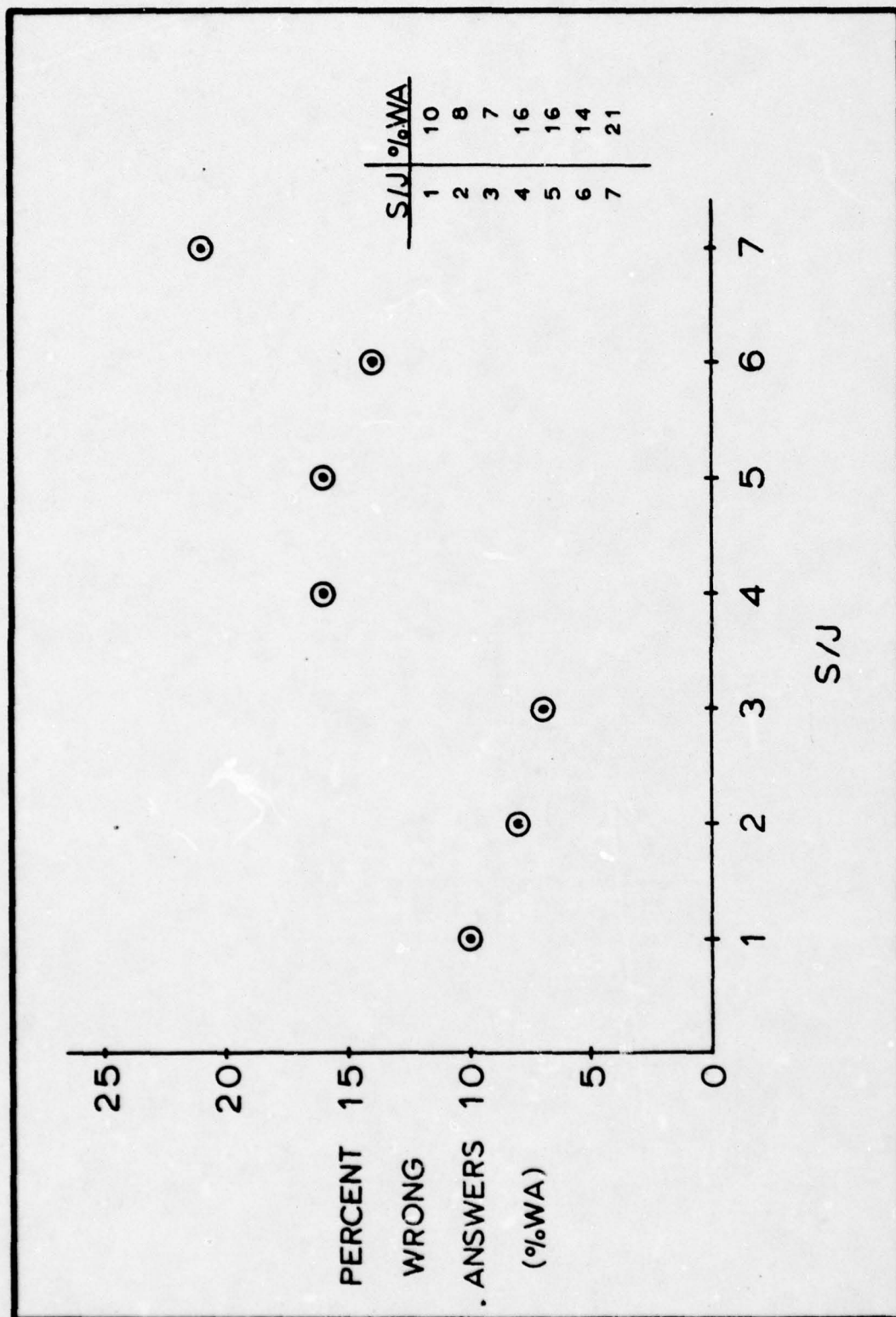


Figure 7. Plot of Percent of Listener Wrong Answers Versus S/J Ratio

Pearson's Correlation Coefficient (P) is used to determine the degree of correlation between the computer calculated mean squared error values and the results of the human listener tests. The formula used to calculate P is

$$P = \frac{\sum_{i=1}^7 (X_i - \bar{X})(Y_i - \bar{Y})}{\left[\sum_{i=1}^7 (X_i - \bar{X})^2 \sum_{j=1}^7 (Y_j - \bar{Y})^2 \right]^{1/2}} \quad (6)$$

where

X_i = Mean Squared Error value

Y_i = error score of human listeners

$$\bar{X} = 1/7 \sum_{i=1}^7 X_i$$

$$\bar{Y} = 1/7 \sum_{i=1}^7 Y_i$$

In this case P was calculated to be 0.74. To find the percentage of the variance in the listener errors accounted for by observing the mean squared error values, under a Gaussian assumption (zero mean Gaussian), P must be squared and multiplied by 100, which gives a value of 55%.

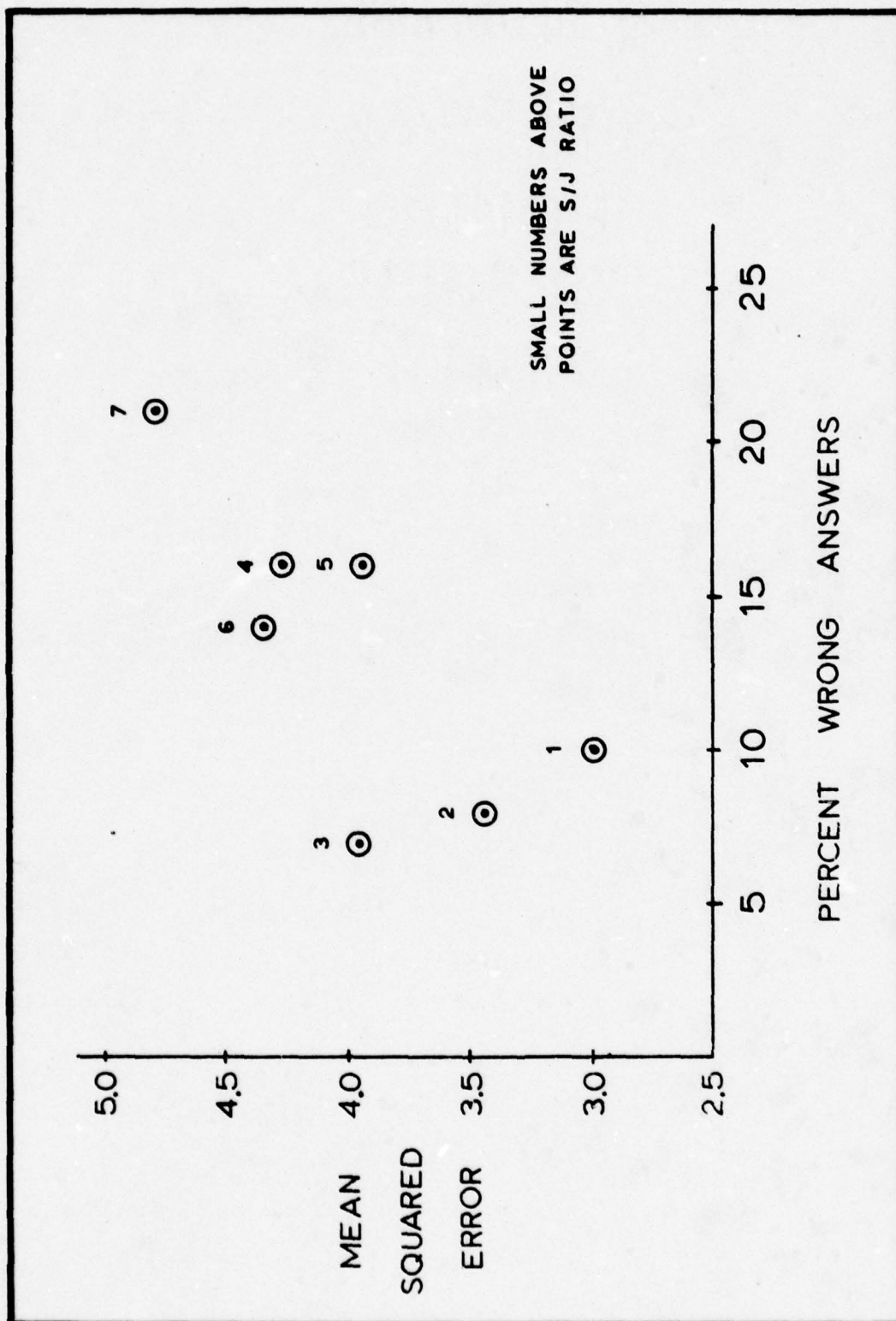


Figure 8. Scatter Plot of Mean Squared Error Versus Percentage of Listener Wrong Answers

VII. Conclusions and Recommendations

The value of 0.74 which was calculated for Pearson's Correlation Coefficient confirms that there is significant correlation between the intelligibility scores of the human listeners and the MSE values calculated by the computer program. This appears to support the ideas set forth in this thesis as a reasonable approach to predicting voice intelligibility. This single comparison is not enough to determine an exact relationship between the mean squared error value and intelligibility, but it seems to provide a starting place for further refinement of the procedure.

The human listener scores used as a comparison were the average of ten listeners. These results showed an unexpectedly high error rate at the best S/J level which decreased as the noise got worse for the first three S/J levels, Figure 7. This trend disappears if a much larger group of listener scores are averaged together and the error rate increases monotonically as the S/J ratio decreases (Ref 2). This listener performance would have given closer agreement with the MSE predictions. This suggests that a listener group considerably larger than ten is required to get a reliable intelligibility figure.

One of the original goals of this thesis was to provide a means of making computerized intelligibility measurement that do not require any special or unique equipment. This goal was met, with the exception of the processing described in section III done by the ASD Analog/Hybrid Branch. A recommendation for further work in this area is to develop a program for one of the more elaborate minicomputers in common use in the Air Force that can perform these operations.

Further investigation in this area could focus on methods of comparing the master word and noisy word, after they have been preprocessed by the ear models, other than MSE. This may yield a better predictor than MSE.

Bibliography

1. Acoustical Society of America. American National Standard Methods for the Calculation of the Articulation Index. ANSI-S3.5-1969. New York: American National Standards Institute, January, 1969.
2. Bauer, John E. Evaluation of Two Sampled Data Communication Systems. M.S. Thesis GE/EE/77-8. Wright-Patterson Air Force Base, Ohio: Air Force Institute of Technology (1977).
3. Dailey, Keith G. and F. Sutton. An Automatic Speech Recognition System Using A Vocoder Input. M.S. Thesis GE/GGC/EE/72-18. Wright-Patterson Air Force Base, Ohio: Air Force Institute of Technology (1972).
4. Electronic Systems Division, U. S. Air Force. Comparative Evaluations of Speech Intelligibility Performance of Three Narrowband Voice Communications Methods: TRIVOC, LPC, and PLPC. ESD-TR-77-131. Hanscom Air Force Base, Massachusetts: ESD, USAF; December, 1976.
5. Hall, William B., Jr. A Digest and Reference Organization of Fast Fourier Transform Literature and Software. VNA, Internal Memo 72-2. Wright-Patterson Air Force Base, Ohio: Analog/Hybrid Systems Branch Computer Center, February, 1972.
6. Hartman, W. J. and S. Boll. Voice Channel Objective Evaluation Using Linear Predictive Coding. FAA-RD-75-189. Washington: Federal Aviation Administration, August, 1976.
7. Hartman, W. J. and K. Gamauf. Objective Evaluation of Voice Communication Channels Using Linear Predictive Coding. unpublished paper. Boulder, Colorado: U. S. Department of Commerce, June, 1976.
8. House, A. S. and G. Hughes. Speech Analysis. AFCRL-69-0371. Bedford, Massachusetts: Air Force Cambridge Research Laboratories August, 1969. AD 696599.
9. Kabrisky, Matthew. A Proposed Model for Visual Information Processing in the Human Brain. Urban, Illinois: University of Illinois Press, 1966.
10. Kiang, Nelson, Y. Discharge Patterns of Single Fibers in the Cat's Auditory Nerve. Cambridge, Massachusetts: The MIT Press, 1965.
11. Laefoged, Peter. Elements of Acoustic Phonetics. Chicago, Illinois: The University of Chicago Press, 1962.
12. Neyman, Ralph W. Computer Recognition of Phonemes in Continuous Speech. M.S. Thesis GE/EE/76-10. Wright-Patterson Air Force Base, Ohio: Air Force Institute of Technology (1976).

13. Niederjohn, R. J. and J. Grotelueschen. "The Enhancement of Speech Intelligibility in High Noise Levels by High-Pass Filtering Followed by Rapid Amplitude Compression." IEEE Transactions on Acoustics, Speech, and Signal Processing, ASSP-24: 277-282 (August, 1976).
14. Ullman, J. R. Pattern Recognition Techniques. New York: Crane, Russak and Co., Inc. 1973.
15. Voiers, W. D.; A. Shappley; and C. Hehmsoth. Research on Diagnostic Evaluation of Speech Intelligibility. AFCRL-72-0694. Bedford Massachusetts: Air Force Cambridge Research Laboratories, January, 1973.

Appendix A

Sequence Chart for Intelligibility Prediction

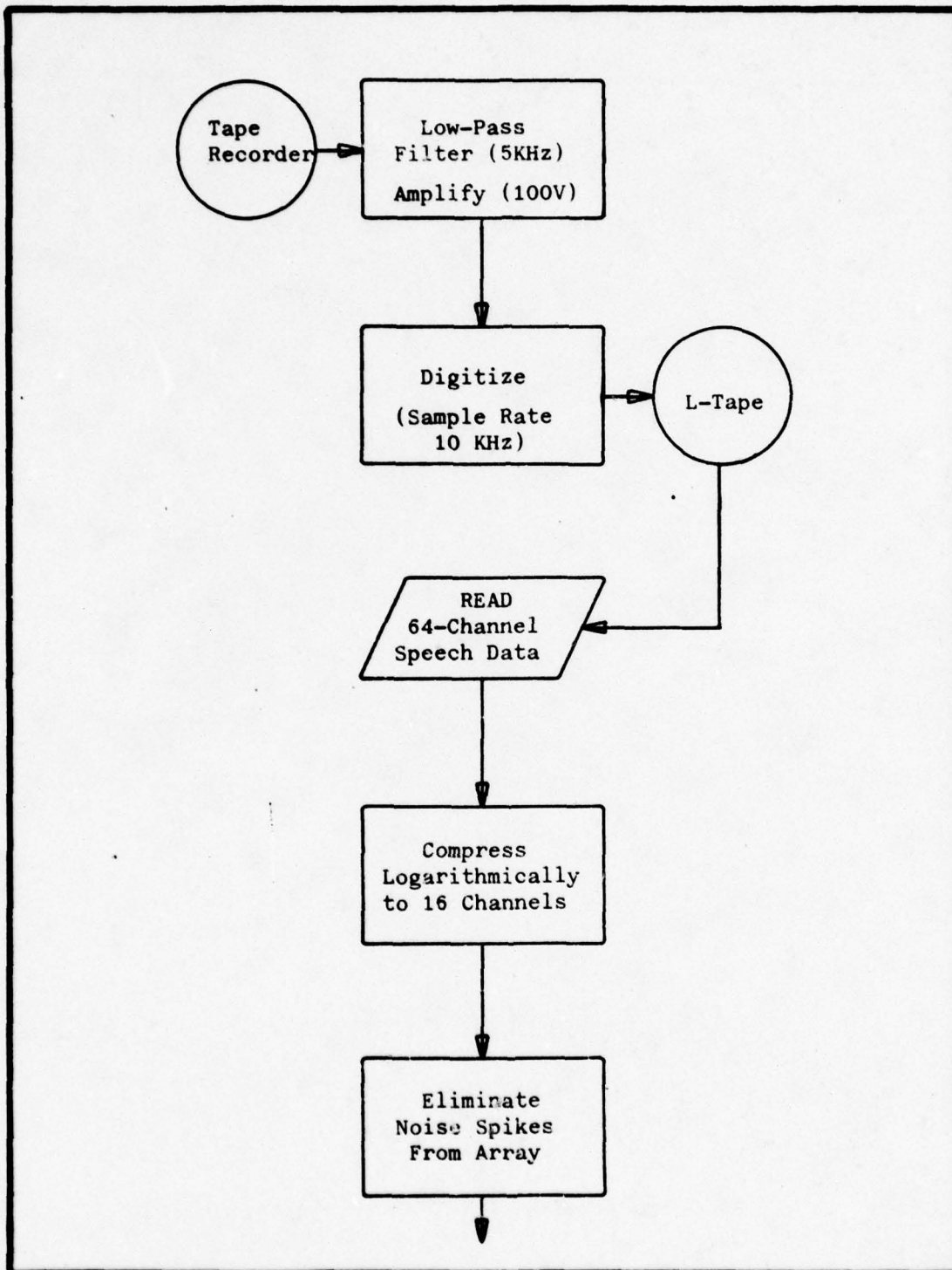


Figure 9. Sequence Chart for Master Tape Processing Program
(Plate 1)

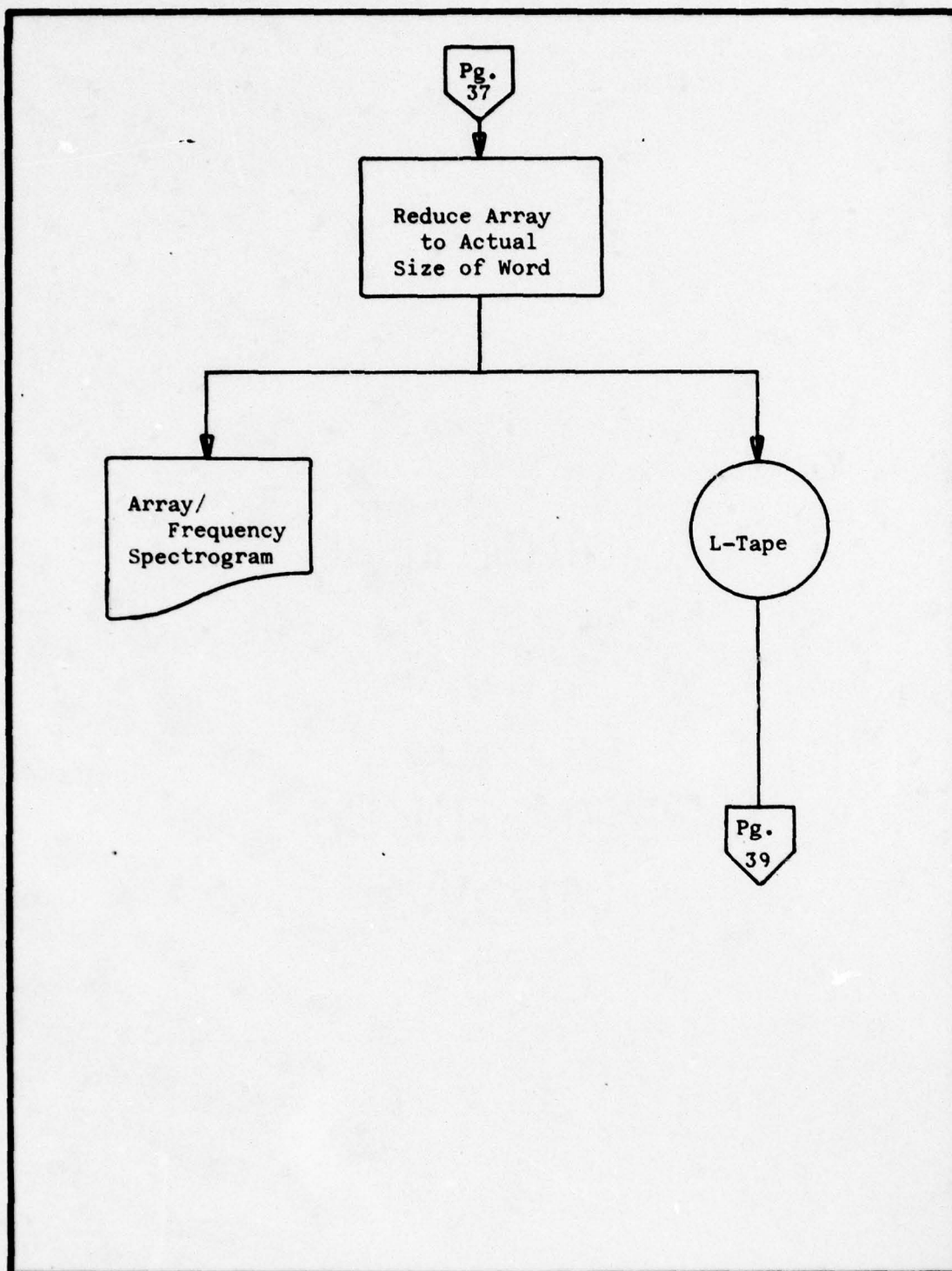


Figure 10. Sequence Chart for Master Tape Processing Program
(Plate 2)

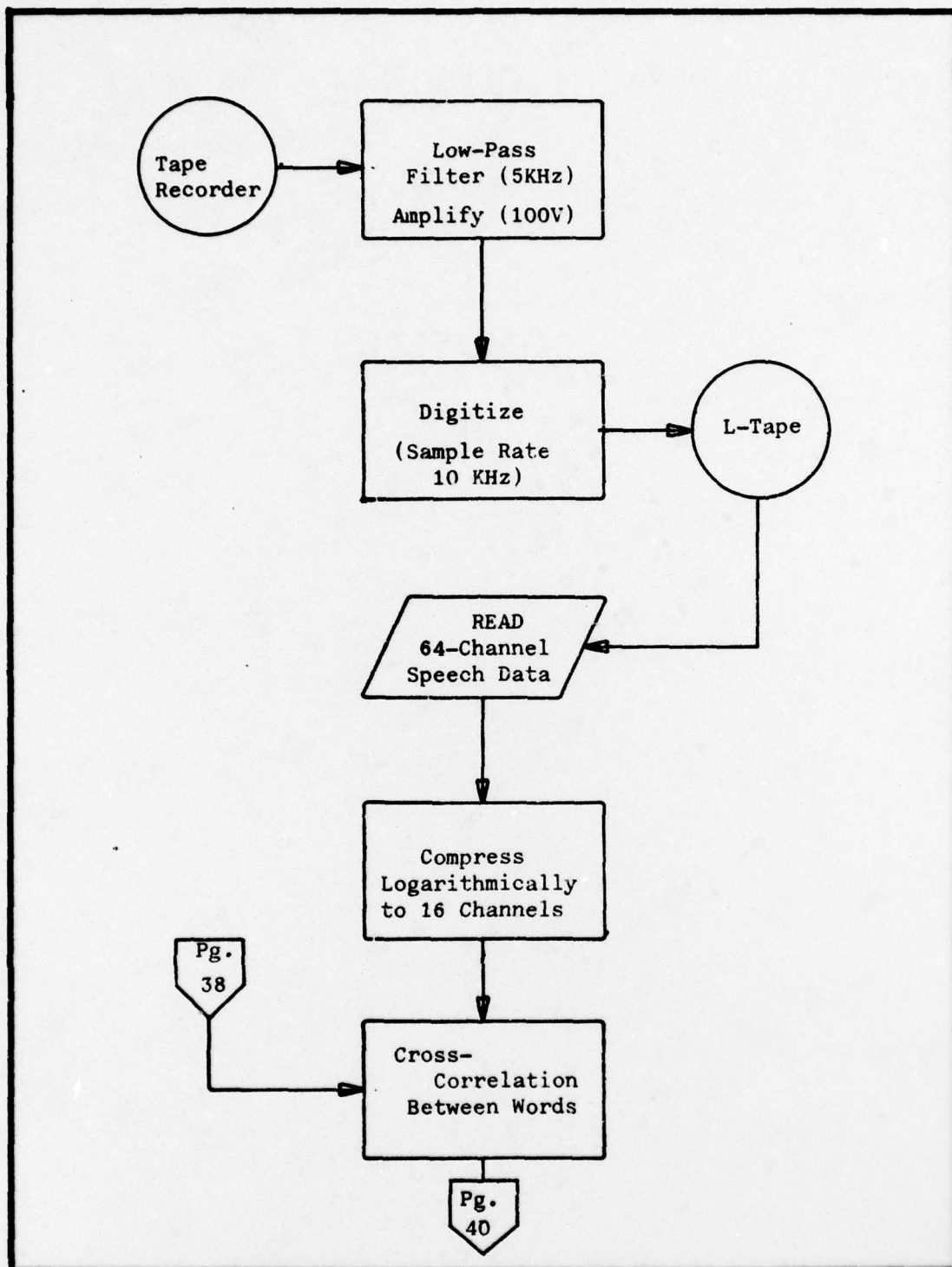


Figure 11. Sequence Chart for Cross-Correlation and MSE Program (Plate 1)

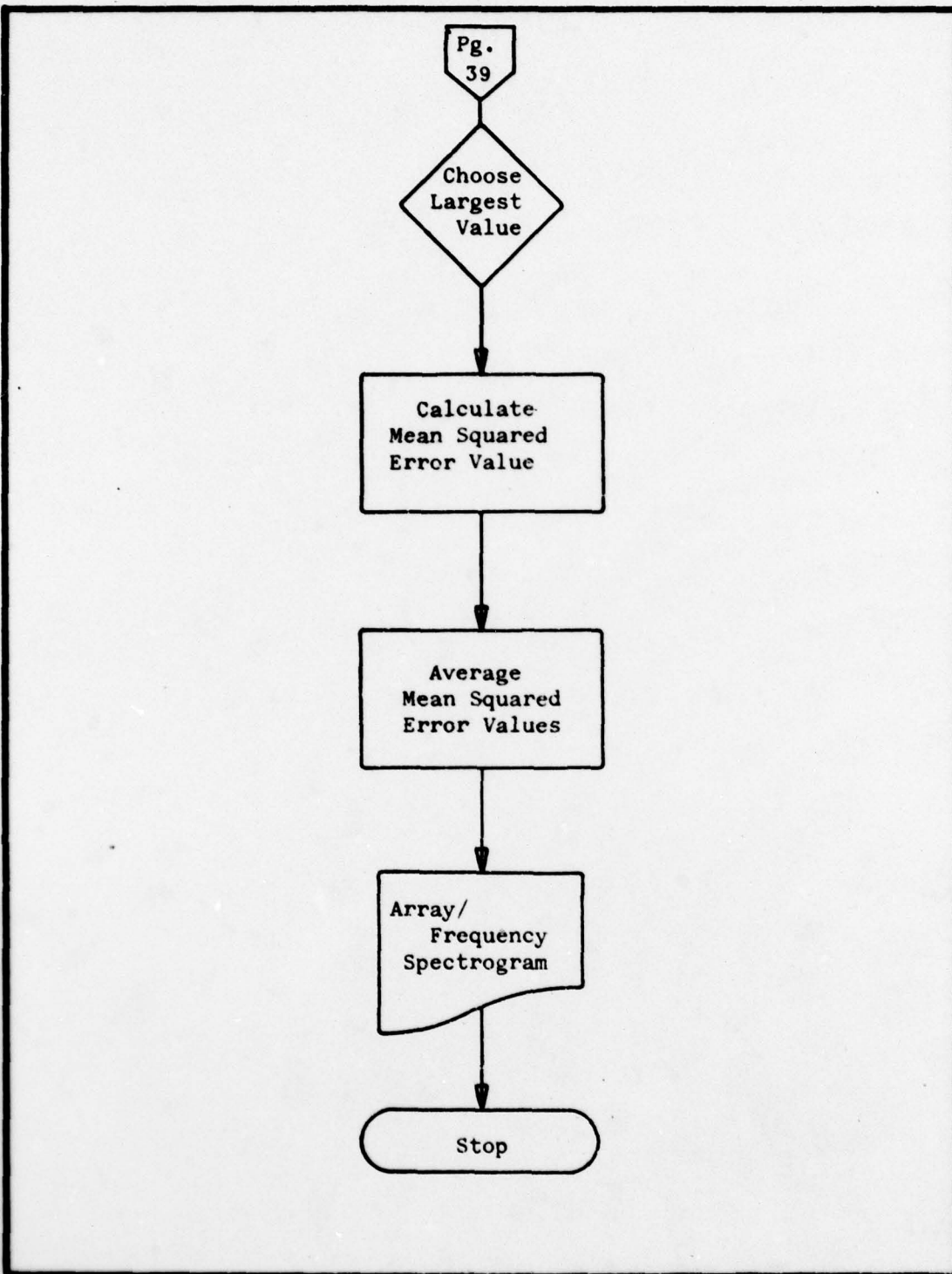


Figure 12. Sequence Chart for Cross-Correlation and MSE Program (Plate 2)

Appendix B

Computer Programs

11/26/77 13.15.40

FIN 4.5414

7-774 OPT=2

```

60 JJ=1
   DO 40 J=1,6
   R(J)=A(J)
   JJ=JJ+1
   CONTINUE
40 DO 50 J=7,11,2
   R(J)=(R(J)+A(J+1))
   JJ=JJ+1
   CONTINUE
50 DO 60 J=13,17,4
   R(J)=(R(J)+A(J+1))+A(J+2)+A(J+3))
   JJ=JJ+1
   CONTINUE
60 DO 70 J=21,25
   R(J)=(R(J)+A(J))
   CONTINUE
70 R(J)=S(J)
   S(J)=S(J)+1
   S(J)=S(J)+1
   DO 80 J=26,31
   S(J)=(S(J)+A(J))
   CONTINUE
80 JJ=JJ+1
   S(J)=S(J)+2
   S(J)=S(J)+2
   DO 90 J=32,40
   S(J)=(S(J)+A(J))
   CONTINUE
90 JJ=JJ+1
   R(J)=S(J)
   S(J)=S(J)+3
   S(J)=S(J)+3
   DO 100 J=41,50
   S(J)=(S(J)+A(J))
   CONTINUE
100 JJ=JJ+1
   R(J)=S(J)
   S(J)=S(J)+4
   S(J)=S(J)+4
   DO 110 J=51,64
   S(J)=(S(J)+A(J))
   CONTINUE
110 JJ=JJ+1
   R(J)=S(J)
500 CONTINUE
C.....
C** NOISE SPIKES ELIMINATED FROM ARRAY.
C.....
120
130 NSUM=0
   DO 400 JJ=1,16
   IF(R(JJ).LT.1.5) R(JJ)=0
   PATR(JJ,1)=R(JJ)
   IF(PATR(JJ,1).GT.0.1) NSUM=NSUM+1
400 CONTINUE
140 NSUM=NSUM
150 CONTINUE

```


11/26/77 03.15.00

FTN 4.5441

PROJ=14 UCTAVE 74/74 OPT=2

```
100  
170  
PRINT 211,(SYMBOLS(IGI(JJ)),JJ=1,16)  
PRINT 211,(SYMBOLS(IGI(JJ)),JJ=1,16)  
FORMAT('m',91X,16A1)  
211 CONTINUE  
212  
213  
214  
215  
216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431  
432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604  
605  
606  
607  
608  
609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695  
696  
697  
698  
699  
700  
701  
702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025  
1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079  
1080  
1081  
1082  
1083  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099  
1100  
1101  
1102  
1103  
1104  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119  
1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133  
1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349  
1350  
1351  
1352  
1353  
1354  
1355  
1356  
1357  
1358  
1359  
1360  
1361  
1362  
1363  
1364  
1365  
1366  
1367  
1368  
1369  
1370  
1371  
1372  
1373  
1374  
1375  
1376  
1377  
1378  
1379  
1380  
1381  
1382  
1383  
1384  
1385  
1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429  
1430  
1431  
1432  
1433  
1434  
1435  
1436  
1437  
1438  
1439  
1440  
1441  
1442  
1443  
1444  
1445  
1446  
1447  
1448  
1449  
1450  
1451  
1452  
1453  
1454  
1455  
1456  
1457  
1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511  
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619  
1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673  
1674  
1675  
1676  
1677  
1678  
1679  
1680  
1681  
1682  
1683  
1684  
1685  
1686  
1687  
1688  
1689  
1690  
1691  
1692  
1693  
1694  
1695  
1696  
1697  
1698  
1699  
1700  
1701  
1702  
1703  
1704  
1705  
1706  
1707  
1708  
1709  
1710  
1711  
1712  
1713  
1714  
1715  
1716  
1717  
1718  
1719  
1720  
1721  
1722  
1723  
1724  
1725  
1726  
1727  
1728  
1729  
1730  
1731  
1732  
1733  
1734  
1735  
1736  
1737  
1738  
1739  
1740  
1741  
1742  
1743  
1744  
1745  
1746  
1747  
1748  
1749  
1750  
1751  
1752  
1753  
1754  
1755  
1756  
1757  
1758  
1759  
1760  
1761  
1762  
1763  
1764  
1765  
1766  
1767  
1768  
1769  
1770  
1771  
1772  
1773  
1774  
1775  
1776  
1777  
1778  
1779  
1780  
1781  
1782  
1783  
1784  
1785  
1786  
1787  
1788  
1789  
1790  
1791  
1792  
1793  
1794  
1795  
1796  
1797  
1798  
1799  
1800  
1801  
1802  
1803  
1804  
1805  
1806  
1807  
1808  
1809  
1810  
1811  
1812  
1813  
1814  
1815  
1816  
1817  
1818  
1819  
1820  
1821  
1822  
1823  
1824  
1825  
1826  
1827  
1828  
1829  
1830  
1831  
1832  
1833  
1834  
1835  
1836  
1837  
1838  
1839  
1840  
1841  
1842  
1843  
1844  
1845  
1846  
1847  
1848  
1849  
1850  
1851  
1852  
1853  
1854  
1855  
1856  
1857  
1858  
1859  
1860  
1861  
1862  
1863  
1864  
1865  
1866  
1867  
1868  
1869  
1870  
1871  
1872  
1873  
1874  
1875  
1876  
1877  
1878  
1879  
1880  
1881  
1882  
1883  
1884  
1885  
1886  
1887  
1888  
1889  
1890  
1891  
1892  
1893  
1894  
1895  
1896  
1897  
1898  
1899  
1900  
1901  
1902  
1903  
1904  
1905  
1906  
1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943  
1944  
1945  
1946  
1947  
1948  
1949  
1950  
1951  
1952  
1953  
1954  
1955  
1956  
1957  
1958  
1959  
1960  
1961  
1962  
1963  
1964  
1965  
1966  
1967  
1968  
1969  
1970  
1971  
1972  
1973  
1974  
1975  
1976  
1977  
1978  
1979  
1980  
1981  
1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997  
1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051  
2052  
2053  
2054  
2055  
2056  
2057  
2058  
2059  
2060  
2061  
2062  
2063  
2064  
2065  
2066  
2067  
2068  
2069  
2070  
2071  
2072  
2073  
2074  
2075  
2076  
2077  
2078  
2079  
2080  
2081  
2082  
2083  
2084  
2085  
2086  
2087  
2088  
2089  
2090  
2091  
2092  
2093  
2094  
2095  
2096  
2097  
2098  
2099  
2100  
2101  
2102  
2103  
2104  
2105  
2106  
2107  
2108  
2109  
2110  
2111  
2112  
2113  
2114  
2115  
2116  
2117  
2118  
2119  
2120  
2121  
2122  
2123  
2124  
2125  
2126  
2127  
2128  
2129  
2130  
2131  
2132  
2133  
2134  
2135  
2136  
2137  
2138  
2139  
2140  
2141  
2142  
2143  
2144  
2145  
2146  
2147  
2148  
2149  
2150  
2151  
2152  
2153  
2154  
2155  
2156  
2157  
2158  
2159  
2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213  
2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267  
2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321  
2322  
2323  
2324  
2325  
2326  
2327  
2328  
2329  
2330  
2331  
2332  
2333  
2334  
2335  
2336  
2337  
2338  
2339  
2340  
2341  
2342  
2343  
2344  
2345  
2346  
2347  
2348  
2349  
2350  
2351  
2352  
2353  
2354  
2355  
2356  
2357  
2358  
2359  
2360  
2361  
2362  
2363  
2364  
2365
```



```

172      PRINT 52
      GO 72 12005,7
      PRINT 53,1,MUSCLE(I)
      GO 73 12005,7
      PRINT 54,1,MUSCLE(I)
      PRINT 55
      PRINT 72
      FORMAT(39A,"****VOICING****",/)
      PRINT 52
      GO 73 12029,7
      PRINT 53,1,MUSCLE(I)
      GO 74 12005,7
      PRINT 54,1,MUSCLE(I)
      PRINT 55
      PRINT 74
      FORMAT(39A,"****ANALITY****",/)
      PRINT 52
      GO 75 12029,7
      PRINT 53,1,MUSCLE(I)
      GO 76 12029,7
      PRINT 54,1,MUSCLE(I)
      PRINT 55
      PRINT 76
      FORMAT(39A,"****SUSTENTION****",/)
      PRINT 52
      GO 77 12005,7
      PRINT 53,1,MUSCLE(I)
      GO 78 12005,7
      PRINT 54,1,MUSCLE(I)
      PRINT 55
      PRINT 78
      PRINT 76
      FORMAT(39A,"****SIBILATION****",/)
      PRINT 52
      GO 79 12005,7
      PRINT 53,1,MUSCLE(I)
      GO 80 12005,7
      PRINT 54,1,MUSCLE(I)
      PRINT 55
      PRINT 80
      PRINT 76
      FORMAT(39A,"****GRAVNESS****",/)
      PRINT 52
      GO 81 12029,7
      PRINT 53,1,MUSCLE(I)
      GO 82 12005,7
      PRINT 54,1,MUSCLE(I)
      PRINT 55
      PRINT 82
      FORMAT(39A,"****COMPACTNESS****",/)
      PRINT 52
      GO 83 12029,7
      PRINT 53,1,MUSCLE(I)
      GO 84 12005,7
      PRINT 54,1,MUSCLE(I)
      PRINT 55
      PRINT 84
      PRINT 82

```

Vita

Wayne R. Beeson was born on 29 October 1943 in Dodge City, Kansas. He attended primary and secondary school in Minneola, Kansas, graduating in 1961. In 1966 he received his undergraduate degree in Chemistry and Mathematics from Northwestern State College, Alva, Oklahoma. After graduation from Northwestern, he entered the Air Force and received his commission through the OTS program on 25 November 1966. From December 1966 until October 1967, he attended the Communications Officer Technical School at Keesler AFB, Mississippi. His first assignment was with the 2063rd Communications Group, Lindsey AS, Wiesbaden, Germany. Captain Beeson spent four years in Germany, serving in both the operations and maintenance branches of the 2063rd. In 1971 he was assigned to Air Force Recruiting Service, Detachment 702, Des Moines, Iowa, as the support services officer. In 1975 he entered the Air Force Institute of Technology in the electrical engineering ES Program.

Permanent Address: RFD
Minneola, Kansas 67865

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER AFIT/GE/EE/77-9	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) AN ALGORITHM FOR DETERMINING SPEECH INTELLIGIBILITY		5. TYPE OF REPORT & PERIOD COVERED MS Thesis
		6. PERFORMING ORG. REPORT NUMBER
7. AUTHOR(s) Wayne R. Beeson, Capt, USAF		8. CONTRACT OR GRANT NUMBER(s)
9. PERFORMING ORGANIZATION NAME AND ADDRESS Air Force Institute of Technology (AFIT/EN) Wright-Patterson AFB, Ohio 45433		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS Project 7071-00-12
11. CONTROLLING OFFICE NAME AND ADDRESS Mr. Charles Jacobs (AFCS/OA) Hq Air Force Communications Service Scott AFB, Illinois 62225		12. REPORT DATE December 1977
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		13. NUMBER OF PAGES 52
		15. SECURITY CLASS. (of this report) Unclassified
		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
16. DISTRIBUTION STATEMENT (of this Report) Approved for public release; distribution unlimited		
17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)		
18. SUPPLEMENTARY NOTES Approved for public release; TAW AFR 190-17 <i>[Signature]</i> JERRAL F. GUESS, Captain, USAF Director of Information		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) Speech Intelligibility Voice Intelligibility Algorithm for Determining Speech Intelligibility		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) A method of predicting speech intelligibility using computer algorithms is presented. Diagnostic Rhyme test number four was used to measure speech intelligibility using a subjective listener test and these results were used as a basis for comparison with the intelligibility predictions made by the computer algorithm. An audio recording of a speaker reading the Diagnostic Rhyme test was made. This recording was run through a General Electric radio system and varying amounts of noise were added. The output of the radio system was re-recorded, providing a copy of the input word corrupted by both additive noise and		

DD FORM 1 JAN 73 1473

EDITION OF 1 NOV 55 IS OBSOLETE

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

radio system distortion effects. Both the input recording and the noisy output recording were digitized by sampling the analog waveforms at a 10 kilohertz rate. These digital samples were converted to a frequency format by windowing the time samples with a rectangular window 128 time samples in length and processing them using Fast Fourier transform techniques. This procedure simulated running the analog speech signal through a bank of contiguous narrow bandpass filters covering the range of 0 to 5 KHz, with center frequencies 78 Hz apart. The output of this process was a matrix array, corresponding to each word from the tape, of amplitude values 200 time windows long and divided into 64 frequency bands. These 64 frequency bands were then combined into 1/3 octave groups to model the frequency sensitivity of the average human ear, which reduced the matrix array to 16 frequency bands. This processing of the analog signal was used to model the preprocessing which occurs in the human ear. A comparison between each word from the input tape and the noisy output tape was then made using a weighted mean squared error calculation. This comparison was conjectured to provide a difference measure which is inversely related to intelligibility. This comparison was used to represent how intelligible the input received from the inner ear is to the brain.

Comparison of the intelligibility results from the human listener tests with the computer processing method outlined above gave a Pearson's Correlation Coefficient value of 0.74 which indicates the computer prediction accounted for 55% of the variance in the listener error scores.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)